# Machine Learning Classification Techniques for Job Titles: A Performance Comparison

Abhishek Nagpal*, Priya Minocha**

*(Department of Information Science and Engineering, Rashtriya Vidyalaya College of Engineering, Bengaluru
Email: abhisheknagpal6633@gmail.com )
**(Department of Information Science and Engineering, Rashtriya Vidyalaya College of Engineering, Bengaluru
Email: priyaminocha24@gmail.com)

----------------------------------\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*----------------------------------

## Abstract:

An unseemly competitor shortlisted and a potential one missed essentially implies an improper resume connected to the erroneous catchphrase. Report arrangement is by and large exorbitantly explored nowadays, because of developing interest in text characterization which has become a significant supporter of the online messages and archives. The redundant errands of an individual arranging the subtleties can be dealt with by the hardware utilizing a specialist framework that accurately catches and recognizes the content and afterward groups it into various classes characterized. After the preprocessing of the information, the characterization is done as a similar examination of Bernoulli's Naïve Bayes, Multinomial Naïve Bayes, Random Forest, Linear SVM and LSVM with flexible punishment arrangement on the Top 30 Job posting dataset with various boundaries and in this way we can break down the conditions between various terms in classes with fluctuating densities and records. The exactness was assessed and LSVM gave the best precision in arranging position privileges dependent on the questions submitted and had the option to accomplish an exactness of approximately 92% for a huge dataset.

*Keywords* —**RF, BNB,MNB,SVM,LSVM.**

----------------------------------\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*----------------------------------

## I. INTRODUCTION

Regular huge numbers of employment forms are submitted. It is preposterous to expect to think about every single application physically. Subsequently organizations use machine learning calculations to order them. All around made requests for employment successfully match work searchers with businesses that are anxious to fill occupations with qualified individuals. AI applications in industry much of the time have the need to handle extremely huge datasets.For most bosses, the primary motivation behind the application is to get evident data about work searchers, like their contact data, qualification, degrees, licenses, abilities, capabilities, business history and a rundown of references. Most importantly, records can be consequently circulated into various spatial areas and can be chronicled quickly with no trade off of human time.In this task, a framework has been carried out for order of reports dependent on text grouping. Record characterization is by and large unnecessarily explored nowadays, because of

developing interest in text grouping which has become a significant supporter of the online messages and reports. It is a programmed way to deal with work title grouping utilizing sets of expectations that decrease the matter of text report arrangement with machine learning . The methodology requires marking the archives with predefined classes to frame a gathering of training information.

Document Classification can be space-explicit and area subordinate information that can be granted to our framework for grouping of particular records just and along these lines, it performs explicit assignments productively. Record grouping can be utilized to separate a ton of heterogeneous archives from the information base to remove various reports and afterward characterize them dependent on the same visual quality [8], tables, designs, and so on. This can improve the report picture recovery master framework and make it more hearty. Report characterization has been broadly concentrated in machine learning, information preparing and information recovery networks and has likewise been applied in different modern settings. A programmed way to deal with work title characterization diminishes to the matter of text report grouping with AI. An AI based archive arrangement approach requires marking records with predefined classes to shape a gathering of training information. A few characterization calculations like Support Vector Machine (SVM) [1] and k-Nearest-Neighbor (kNN) [2] have been applied to various modern issues with great observational outcomes. AI applications in industry regularly have the need to handle enormous datasets. A few huge scope appropriated AI systems have adjusted machine learning calculations to circulate and increase structures [3]. These frameworks are appropriate to the unquestionably fundamental application circumstances that incorporate immense datasets and have implied computationally mentioning readiness of complex learning models in a way as close to ongoing derivation limitations. In this paper, a framework

has been proposed of order of records dependent on text grouping. They expect to produce the content vector utilizing TF-IDF vector strategy in which the term recurrence is development of a word reference of words that are happening in the report and afterward building a vector out of it. This vector is then likewise developed for backwards recurrence words, which implies the numeric worth of the reports where the term can happen. In light of these two vectors which are developed out of the extricated highlights, preprocessing them and changing them to an ideal structure learning calculations are applied to them, which are Bernoulli's Naïve Bayes (BNB), Multinomial Naïve Bayes (MNB), Random Forest (RF) and Linear SVM (LSVM) for order into some classes of the work titles given to the workers.Chi-Square Test is utilized for all out highlights in a dataset. It decides whether the relationship between two absolute factors of the example would mirror their genuine relationship in the populace.

### A. Motivation

In recent times the management of the job applications was becoming unmanageable since it is difficult to keep a track of availability at a particular place and time as the requirement varies from company to company and from place to place. Everyday there are thousands of job applications are submitted from which many times the deserving candidates are not recognized. There has been a need for an advanced machine learning classification algorithm that can predict the job title to some extent by the job description.

## II. LITERATURE SURVEY

The following papers were surveyed to get a better understanding of the topic-

Ms. Yoginee R. Surkar and Prof. S. W. Mohod in their paper , "A Review on Feature Selection and Document Classification utilizing Support Vector Machine"proposed a book grouping dependent on the element choice and preprocessing by lessening the dimensionality of the element vector and expanding the order exactness. This has additionally been applied to many example grouping issues, for example, picture acknowledgment, discourse acknowledgment, text arrangement, face identification and flawed card discovery, and so forth Example acknowledgment intends to order information dependent on either deduced information or measurable data removed from crude information, which is an amazing asset in information partition in numerous disciplines[1].

Bei Yu and Linchi Kwok in their paper, " Classifying Business Marketing Messages on Facebook", proposed that the SVM classifier is prepared to naturally isolate these two sorts of messages, expecting to utilize this instrument to investigate messages from numerous organizations and subsequently screen the advancement of their online media use over time.It is tracked down that the classifier prepared with tf-idf weighted grammatical form highlights performed best. It is superior to classifiers prepared with word features[2].

Dr. Shaukat Wasi and Dr. Zubair Ahmed Shaikh in their paper, "Context‐based email arrangement model" , the analysts fostered a setting based occasion email classification model that classifies messages into specific envelopes. Such classification is significantly attractive for supporting occasion arranging and the executives errands. Experimentation results show that utilizing layout based methodology upheld by chart coordinating and mining calculations creates high classification exactness for the occasion type based email classification [3].

Gerard Salton and Christopher Buckley in their paper, "Term-weighting approaches in programmed text recovery", showed an exploratory proof in the course of recent years demonstrates that text ordering framework dependent on the task of the suitably weighted single terms produce recovery results that are better than those possible with more elaborate content portrayals. These outcomes rely critically upon the decision of a powerful term weighting framework. The paper sums up the experiences acquired in programmed term weighting and gives benchmark single term ordering models [4].

Andrew McCallum and Kamal Nigam in their paper, "A Comparisons of Event Models for Naive Bayes Text Classification " looked at the hypothesis and practice of first request probabilistic classifiers, the two of which make the Naive Bayes presumption . The multinomial model is discovered to be consistently better compared to themulti-variate Bernoulli model [5].

Anjali G Jivani in her paper, "A Comparative Study of Stemming Algorithms", talked about stemming and its advantages. Numerous calculations were talked about and they showed a ton of likenesses. The principle contrast lies in utilizing either a standard based methodology or an etymological one. A factual stemmer might be a language autonomous yet doesn't generally give a solid and right stem. The issue of over stemming and under stemming can be decreased just if sentence structure just as semantics of the words and their position is contemplated [6].

## III. METHODOLOGY

The stepwise procedure is mentioned below-

### A. Articulate the issue early

Understanding what you need to foresee will assist you with choosing which information might be more important to gather. While forming the issue, lead information investigation and attempt to think in the classifications of arrangement, bunching, relapse, and positioning that we discussed in our whitepaper on business use of machine learning. In layman's terms, these undertakings are separated in the accompanying way:

**Classification-** You need a calculation to answer parallel yes-or-no inquiries or you need to make a multiclass grouping . You likewise need the correct answers named, so a calculation can gain from them. Check our guide on the most proficient method to handle information marking in an association.

**Clustering-** You need a calculation to discover the principles of characterization and the quantity of classes. The primary contrast from arrangement assignments is that you don't really have the foggiest idea what the gatherings and the standards of their division are. For example, this generally happens when you need to section your clients and tailor a particular way to deal with each portion contingent upon its characteristics.

**Regression-** You need a calculation to yield some numeric worth. For instance, in the event that you invest an excessive amount of energy thinking of the correct cost for your item since it relies upon numerous components, relapse calculations can help in assessing this worth.

**Ranking** - Some machine learning calculations simply rank articles by various highlights. Positioning is effectively used to suggest motion

pictures in video web-based features or show the items that a client may buy with a high likelihood dependent on their past search and buy exercises.

### B. Set up information assortment components

Making an information driven culture in an association is maybe the hardest piece of the whole drive. We momentarily shrouded this point in our story on AI procedure. On the off chance that you expect to utilize ML for prescient investigation, the primary thing to do is battle information discontinuity.

### C. Check the information quality

Indeed, even the most complex machine learning calculations can't work with helpless information On the off chance that if the information is gathered or marked by people, check a subset of information and gauge how regularly botches happen.If your named dataset has 1,500 passages named as dependable and just 30 that can be considered untrustworthy, the model will not have enough examples to find out about the questionable ones.

### D. Arrangement information to make it reliable

Information organizing is now and then alluded to as the record design you're utilizing. What's more, this isn't a very remarkable issue to change over a dataset into a record design that accommodates your algorithm framework best.

### E. Reduce the data - Decrease information

In case you're setting up a dataset considering specific assignments, it's smarter to lessen information. Since you understand what the objective trait is, sound judgment will control you further. You can accept which esteems are basic and which will add more measurements and intricacy to your dataset with no gauging commitment. This methodology is called trait testing.

### F. Complete information cleaning

Since missing qualities can substantially decrease expectation exactness, focus on this issue. As far as AI, expected or approximated values are "all the more right" for a calculation than simply missing ones. Regardless of what you don't have the foggiest idea about the specific worth techniques exist to better "expect" which worth is missing or sidestep the issue. Picking up the correct methodology additionally vigorously relies on the information and ideas you have.

Substitute the missing mathematical qualities with mean figures For straight out qualities, you can likewise utilize the most incessant things to fill in.

### G. Decompose the data

A few qualities in your informational index can be mind boggling and disintegrating them into numerous parts will help in catching more explicit connections. This interaction is really the inverse to diminishing information as you need to add new properties dependent on the current ones.

### H. Join value-based and trait information

You may have a few information sources or logs where these kinds of information dwell. The two sorts can improve each other to accomplish more noteworthy prescient force. For example, in case you're following hardware sensor readings to empower prescient support, probably you're producing logs of conditional information, yet you can add such characteristics as the gear model, the clump, or its area to search for conditions between gear conduct and its credits.

### I. Rescale information

Information rescaling has a place with a gathering of information standardization methods that target improving the nature of a dataset by lessening measurements and staying away from the circumstance when a portion of the qualities overweight others. min-max standardization can be utilized. It involves changing mathematical qualities to ranges, e.g., from 0.0 to 1.0 where 0.0 addresses the negligible and 1.0 the greatest qualities to try and out the heaviness of the value characteristic with different traits in a dataset.

### J. Discretize information

At times you can be more successful in your expectations on the off chance that you transform mathematical qualities into absolute qualities. This can be accomplished, for instance, by isolating the whole scope of qualities into various gatherings.

## IV. PERFORMANCE EVALUATION

The proficiency of any machine learning calculation is resolved utilizing measures like genuine positive rate, bogus positive rate, genuine negative rate and bogus negative rate. Affectability and explicitness is utilized to clarify clinical symptomatic tests and gauge how great the test was. Preparing and assessing factual execution on similar information yields over idealistic outcomes. Cross approval improves

execution exactness. Cross approval is a factual strategy that analyzes AI conspires by separating the information into test and train sets. The train set is utilized to prepare the information and the test set is utilized to approve the model.

In k-overlay cross approval , information is parceled into k equivalent size folds. The k-cycles are dealt with and approved to such an extent that inside every emphasis, an alternate crease of the information is waited for approval and the excess k-1 overlap is utilized for learning.

For precision, affectability and particularity , the terms TP,TN,FP and FN are utilized.

- TP is a result where the model accurately predicts the positive class.
- TN is an outcome where the model adequately predicts the negative class.
- FP is a result where the model mistakenly predicts the positive class.
- FN is an outcome in which the model forecasts the negative class wrongly.

### A. Confusion Matrix

A Confusion Matrix is a N X N network utilized for assessing the exhibition of an arrangement model ,where N is the quantity of target classes. The framework contrasts the genuine focus on qualities and those anticipated by the machine learning calculations.



Fig 1. Confusion Matrix Terminology

### B. Requirements for Performance Measure

a. **Accuracy -** Precision is the proportion between True Positives and all Positives. It gives us a proportion of pertinent information**.**

b. **Recall -** The recall is the proportion of our model accurately distinguishing True Positives. Review additionally gives us the proportion of how precisely our model is recognizing the significant information.

c. **F-Score-**F-Score or F1 Score gives a model's precision on a dataset. It is the consonant mean of exactness and review.

d. **Precision -** Accuracy of machine learning characterization calculation is one approach to quantify how regularly the calculation orders an information point accurately. Exactness is the quantity of accurately anticipated information that calls attention to all the information focused**.**

e. **Specificity -** Specificity otherwise called positive rate which is characterized as the extent of real negatives , which got anticipated as negative.

## V. IMPLEMENTATION DETAILS

Several Algorithms were tested on the dataset in order to achieve better performance.

- **TF-IDF -** This represents Term Frequency-Inverse Document Frequency. TF-IDF is a factual measure that lifts how important a word is to an archive in an assortment of reports. This is finished by duplicating two measurements , how frequently a word shows up in a record, and the opposite report recurrence of the word across a bunch of archives.

- **Bernoulli Naive Bayes -** Naive Bayes classifier is a probabilistic classifier which implies that given an info, it gives the likelihood of the information being arranged for every one of the classes. It is otherwise called Conditional Probability.

- **Multinomial Naive Bayes -** Multinomial Naive Bayes classifier is reasonable for arrangement with discrete **Multinomial Naive Bayes -** Multinomial Naive Bayes classifier is reasonable for arrangement with discrete highlights (for example word means text grouping).The multinomial dissemination regularly requires number element tallies.

- **Random Forests -**Random Forest Classifiers are learning techniques for order, relapse and different undertakings that work by developing a huge number of choice trees at preparing time

and yielding the class that is the method of the arrangement.

- **Linear SVM** - Support Vector Machine is a straight model for arrangement and relapse issues. A SVM model is essentially a portrayal of various classes in a hyperplane in multidimensional space. The objective of SVM is to partition the datasets into classes to track down a greatest minimal hyperplane(MMH).

- **K-Nearest Neighbor** -The fundamental thought here is to decide the closest neighbor that suggests the closest plot point in the archive space. A vector is developed for each test set and centroid vector for each class is developed that computes the similitude between each class and the vector developed and afterward characterizes it dependent on the class it coordinates the most or for the class for which similitude is most extreme.

| S.No. | Label | Test | Train | Total | % of Total |
|---|---|---|---|---|---|
| 1. | Project Manager | 239 | 958 | 1197 | 80 |
| 2. | Python Developer | 220 | 871 | 1901 | 80 |
| 3. | Receptionist | 28 | 103 | 131 | 80 |
| 4. | Clerk | 91 | 347 | 438 | 80 |
| 5. | Probationary Officer | 193 | 759 | 952 | 80 |
| 6. | Account Executive | 393 | 1577 | 1970 | 80 |
| 7 | Upkeep Technician | 411 | 1636 | 2047 | 80 |

## VI. RESULTS

There were some algorithms that were tried on our dataset. The dataset was splitted into testing and training data. Stop words are the most well-known words in any normal language. To dissect the content information and build NLP models, these stop words probably won't enhance the significance of the record. By and large , the most usually utilized stop words are "this", "is", "in", "for", "when", "to", "at" etc. To improve our dataset, stemming and lemmatization was performed. Stemming is the way toward lessening a word to its promise stem that joins to postfixes and prefixes or to the underlying foundations of the words to their foundations of the words known as lemma. Stemming is additionally a piece of inquiries and Internet web indexes.

Lemmatization is the way toward changing a word over to its base structure. The contrast among stemming and lemmatization will be, lemmatization considers the specific situation and converts the word to its significant base structure, though stemming simply eliminates the last couple of characters, frequently prompting mistaken implications and spelling blunders.

By using TF-IDF , the features were reduced to 40,000 which is approximately 1/8th of the total sample.

The dataset was cut into preparing and testing dataset in the proportion 80:20.

**TABLE 1**
ILLUSTRATION OF DEFINED HOLD OU INFORMATION PARTING

TABLE 2
RESULTS WITH 5500 SAMPLES

| S.No. | Techniques | Accuracy Scores |
|---|---|---|
| 1. | BNB | 0.412 |
| 2. | MNB | 0.464 |
| 3. | RF | 0.882 |
| 4. | LVSM | 0.920 |
| 5. | LVSM with elastic penalty | 0.911 |



Fig 2. Different Classifiers vs Accuracy Scores

## VII. CONCLUSION

This paper fundamentally examines arrangement calculations dependent on a dataset that groups work titles dependent on the question portrayal. The exploratory outcomes show that the strategy not just in

huge datasets, set palatable order results yet has additionally shown a decent arrangement execution on a little example. Nonetheless, customary Bayesian order calculation dependent on word recurrence in the two diverse datasets has an alternate arrangement impact. The best outcomes on various information sizes are best shown by the straight SVM. While the Naïve Bayes based classifiers inadequately recognized the records, there was anyway a lesser computation and execution time that was used by them when contrasted with other learning procedures. The direct SVM with versatile punishment accomplished correctnesses over 90% even after expanding the quantity of tests which with straight SVM, the best outcome on this dataset could be accomplished with approximately 95% or more of the occasions upon accurately grouping the work privileges. This order framework can be reached out to all the pursuit of employment motors and sites wherein numerous candidates present their subtleties for finding a new line of work. This framework will be useful in effectively characterizing them into classifications that they need to apply for and get coordinates with the most fitting managers. This way CV based report order can control joblessness and assist individuals with landing their ideal positions. Accordingly, there is a great deal of extent of improving the productivity of the current frameworks.

## VIII. FUTURE ENHANCEMENT

This arrangement framework can be stretched out to all pursuit of employment motors and sites wherein numerous candidates present their subtleties for finding a new line of work. This framework will be useful in effectively characterizing them into classes that they need to apply for and get coordinates with the most appropriate managers. The framework will be useful in effectively grouping them into various classifications. It has scope in improving the proficiency of the current framework.

## ACKNOWLEDGEMENT

## REFERENCES

[1] Jehad Ali, Rehanullah Khan, Nasir Ahmad, and Imran Maqsood (2012). Random Forests and Decision Trees are two types of decision trees. IJCSI (International Journal of Computer Science Issues) is a journal that publishes articles on computer science issues.

[2 ] https://www.kaggle.com/bman93/dataset Bman93 - https://www.kaggle.com/bman93/dataset Job Description & Job Title Set of data (April,2020)

[3] B. Scholkopf, S. Mika, C. J. Burges, P. Knirsch, K. R. Muller, G. Ratsch, and A. J. Smola (1999). rf versus input space

[4] R. Jensen, Q. Shen, R. Jensen, R. Jensen, R. Jensen, R. Jen (2004). Dimensionality reduction with semantics preserved: rough and fuzzy-rough techniques 16(12), 1457-1471, IEEE Transactions on Knowledge and Data Engineering.

[5] C. Shin, D. Doermann, and A. Rosenfeld (2001). Structure-based features are used to classify document pages. 3(4), 232-247, International Journal on Document Analysis and Recognition.

[6] N. Chen, D. Blostein, N. Chen, N. Chen, N. Chen, N. Chen, N. Chen, N (2007). A look into document image classification: the challenge, the architecture of the classifier, and the results. 10(1), 1-16, International Journal of Document Analysis and Recognition (IJDAR)..

[7] Xin Jin, Anbang Xu, Rongfang Bie, and Ping Guo (2006). SAGE Gene Expression Profiles: Machine Learning Techniques and Chi-Square Feature Selection for Cancer Classification 106-115. 10.1007/11691730 11 Lect Notes Comput Sci. 3916.

[8] B. Yu and L. Kwok, Yu and Kwok, Yu and Kwok, Yu and Kwok, Yu and (2011). Using Facebook to categorise company marketing communications. Association for Computing Machinery Special Interest Group on Information Retrieval, Beijing, China.

[9] K.Nigam, A McCallum, A McCallum, A McCallum, A McCallum, A McCal For naive Bayes text classification, a comparison of event models is made. Workshop on Learning for Text Categorization at AAAA-98, 2004

[10] G. Salton and C. Buckley (1988). In automatic text retrieval, term-weighting algorithms are used. 513-523 in Information Processing & Management, vol. 24, no. 5.

[11] T. Brückner, P. Suda, H. U. Block, and G. Maderlechner (1996, April). Automatic address and content interpretation for in-house letter distribution. Las Vegas, USA, Proceedings of the 5th Annual Symposium on Document Analysis and Information Retrieval (pp. 67-75)

[12]. J. Pennington, R. Socher, and C. D. Manning (2014, October). Glove: Word representation using global vectors. In Proceedings of the 2014 Empirical Methods in Natural Language Processing (EMNLP) Conference (pp. 1532-1543)

[13] Sofia Visa, Brian Ramsay, Anca Ralescu, and Esther Knaap (2011). Feature Selection Using a Confusion Matrix CEUR Workshop Proceedings, Volume 710, Pages 120-127.

[14] K. Gayathri and A. Marimuthu, Gayathri and Marimuthu, Gayathri and Marimuthu, Gayathri and Marimuth (2013, January). Pre-processing of text documents with the KNN in preparation for classification with the SVM. The seventh International Conference on Intelligent Systems and Control (ISCO) took place in 2013. (pp. 453-457). IEEE is an acronym for "Institute of Electrical and

[15] Y. Goldberg and O. Levy, Y. Goldberg and O. Levy, Y. Goldberg and O. Levy, Y. Gold (2014). word2vec is a programme that converts words into vectors Mikolov et alnegative-sampling .'s word-embedding approach is deduced. arXiv:1402.3722 is an arXiv preprint.

[16] M. J. Meena and K. R. Chandran (2009, December). Positive attributes selected via statistical approach for Naive Bayes text classification. The First International Conference on Advanced Computing was held in 2009. (pp. 28-33). IEEE (Institute of Electrical and Electronics

[17] Random Forests, Machine Learning 45(1), 5-32, 2001.Breiman, L., Random Forests, Machine Learning 45(1), 5-32, 2001.

[18] Jaideepsinh K. and Jatinderkumar Saini. the year (2016) Stop-Word Removal Algorithm and Sanskrit Language Implementation 150. 15-17. 10.5120/ijca2016911462. International Journal of Computer

Applications. 150. 15-17. 10.5120/ijca2016911462.

[19] Quan Zhou, Wenlin Chen, Shiji Song, Jacob Gardner, and Kilian Weinberger. In 2014, With an Application to GPU Computing, a Reduction of the Elastic Net to Support Vector Machines

[20]J. R. Brzezinski and G. J. Knafl (1999, September). Context-based classification using logistic regression modelling. Proceedings are available online. The tenth International Workshop on Database and Expert Systems Applications is a conference on database and expert systems applications. DEXA 99 DEXA 99 DEXA 99 DEXA  (pp. 755-759). IEEE is an acronym for "Institute of Electrical and

[21] J. Pennington, R. Socher, and C. D. Manning (2014, October). Glove: Word representation using global vectors. In Proceedings of the 2014 Empirical Methods in Natural Language Processing (EMNLP) Conference (pp. 1532-1543).

[22] Cesarini, F., Lastri, M., Marinai, S., & Soda, G. Cesarini, F., Lastri, M., Marinai, S., & Soda, G. (2001, September). Modified XY trees are encoded for document classification. The Sixth International Conference on Document Analysis and Recognition's Proceedings (pp. 1131-1136). IEEE is an acronym for "Institute of Electrical and

[23]G. Aghila (2010). A Review of the Naive Bayes  Machine Learning Approach for Text Document Classification preprint arXiv:1003.1795, arXiv:1003.1795, arXiv:1003.1795, arXiv:1003.179

[24] P. Q. Liu and J. J. Feng (2010). The Naive Bayes Text Classification Algorithms have been improved. 187-188 in Microcomputer Information, 26(27).